

Synchronization of Temporal Constructs in Distributed Multimedia Systems with Controlled Accuracy

Sang H. Son and Nipun Agarwal

Department of Computer Science
University of Virginia
Charlottesville, VA 22903
son@virginia.edu, nipun@virginia.edu

Abstract. *With the inception of technology in communication networks such as ATM, it will be possible to run multimedia applications on future integrated networks. Synchronization of the related media data is one of the key characteristics of a multimedia system. In this paper we present a scheme for synchronization of multimedia data across a network where the accuracy of detecting asynchronization and predicting the future asynchrony is variable and can be tailored to the intended application. The multimedia data when sent across the network may also be stored at an intermediate node and later retrieved for display. We extend the scheme and present a mechanism wherein synchronization of all the possible temporal constructs is supported and not restricted to the "in-parallel" construct which is only one of the thirteen possible temporal relations.*

1. Introduction

The innovation that multimedia systems provide is the *integrated* manipulation of multimedia information; the independent manipulation of each separate information medium does not provide anything new. This ability of integrated manipulation of multimedia information provides the potential for dynamically creating new, user-directed composite data sets. The challenge imposed by multimedia is not simply providing I/O support for the various media as audio, video and text. The challenge is of providing support for synchronizing otherwise autonomous data transfers within and across the computers [3]. The support for multimedia service is guided by two distinguishing features, namely support for continuous media and synchronization among media streams [2].

The exploration of the issues of synchronization in a multimedia system and their solutions is in an incipient stage. Steinmetz [4] and Little and Ghafoor [5] have discussed methods for formally describing the synchronization

requirements in a multimedia environment. Steinmetz discussed the characteristics of a multimedia system and presented a set of constructs for expressing inter-media relationships; Little and Ghafoor proposed a strategy for formal specification and modeling of multimedia composition with respect to inter-media timing, based on the logic of timed Petri nets. Hoepner [6] explored the synchronization of multimedia objects for presentation based on the Petri net model. Anderson and Homsy [7] described algorithms for synchronization among interrupt-driven media I/O devices, but they are only applicable to single site multimedia workstations. Nicolaou [8] attempted the synchronization problem in a two-level scheme; by defining explicit synchronization properties at the presentation level and by providing control and synchronization operations at the physical level. Escobar et. al [9] presented an adaptive flow synchronization protocol that permits synchronizing in a distributed environment, but however operates under the assumption that global synchronized clocks are present; which is not a very practical assumption [2] considering that ATM networks will not provide a global synchronized clock. Rangan et. al [2] proposed a feedback technique to detect asynchronization among the media streams in an distributed environment and to steer them to synchrony thereafter. The accuracy with which the multimedia server can detect the instant of playback of a media unit at the destination is however bounded by the network jitter. For wide area network applications this value may be greater than the acceptable asynchrony among the media streams, thereby restricting the purview of the mechanism. For instance, if the maximum permissible asynchrony among the media streams for an application is ζ and the network jitter is d , and if $d > \zeta$, then this particular application cannot be sent over the network.

Furthermore in the approaches of both Escobar and Rangan, the synchronization discussed is primarily synchronization "in-parallel" among the various media streams. By synchronization "in-parallel", we mean the capability to display the various media units simultaneously. Hence if the requirement of an application is that a segment of video

This work was supported in part by ONR, IBM, DOE and CIT.

Report Documentation Page			Form Approved OMB No. 0704-0188		
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE 1994		2. REPORT TYPE		3. DATES COVERED 00-00-1994 to 00-00-1994	
4. TITLE AND SUBTITLE Synchronization of Temporal Constructs in Distributed Multimedia Systems with Controlled Accuracy			5a. CONTRACT NUMBER		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S)			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) University of Virginia, Department of Computer Science, 151 Engineer's Way, Charlottesville, VA, 22094-4740			8. PERFORMING ORGANIZATION REPORT NUMBER		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)			10. SPONSOR/MONITOR'S ACRONYM(S)		
			11. SPONSOR/MONITOR'S REPORT NUMBER(S)		
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES 6	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

is to be played K time units *after* an audio clip - without any constraints on when the audio clip is to commence or when the video clip should end, this requirement cannot be easily accomplished by the synchronization mechanisms which support only the parallel or equal temporal construct.

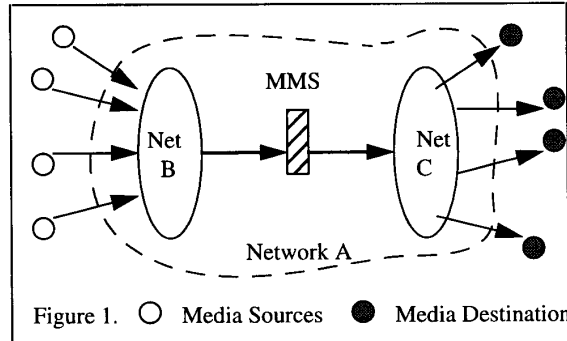
In this paper, we present a scheme for the synchronization of various media streams based on the synchronization marker concept which does not assume the presence of globally synchronized clock. The accuracy of detecting asynchrony in our approach is variable and independent of the network jitter and can thus be tailored to the intended application. We address the issue of assigning relative time stamps to the data originating from different media sources and discuss how our algorithm assigns the time-stamps with an accuracy that can be varied to suit the application. We then present a mechanism that guarantees the synchronization of any possible temporal relation among various media streams, not just the "in-parallel" temporal construct.

2. System Model

Consider a simple scenario, where media data such as audio, video and text is being generated by a number of different media sources. This data is desired to be transmitted across the ATM network to a destination where it is to be played back. The different media data may be routed through different parts of the network which may experience different quality of service, such as jitter and delay. The devices at the destination which play back the media data, may also have mismatches in the playback rates and may differ from the devices used at the source. Furthermore, the data may be stored on a physical storage at an intermediate site and may later be retrieved for play back. In any of the above cases, it is required that the data must be synchronized while it is being played at the destination devices - i.e. played back in the same time order as was generated at the source. A different scenario having similar synchronization requirement would be when data being sent from a single source is required to be received at multiple destinations simultaneously, e.g. in a teleconferencing application.

The above scenario is depicted in Figure 1. The goal is to send some multimedia data from the media sources to the media destinations over the network A. Let a node on the network, called the Multimedia server (MMS) be responsible for controlling the synchronization of media data at the destination and for storing the data if required by the application. The network A is logically divided into network B and network C. The data travels from network B to the MMS and from the MMS travels over network C to

the media destinations. The MMS may hence be at any point on the network. Since the accuracy of detecting the asynchrony of our algorithm is independent of the network jitter, the MMS need not be necessarily physically close to the media destinations.



3. Assigning the Relative Time-stamps

The aim of assigning relative time-stamps (RTS) is that media units in different data streams that have the same relative time-stamp be synchronized during play back. Assigning RTS to the data can be done quite conveniently by counting the number of bytes in each data stream, if compression is not involved. For instance, if a RTS is desired to be assigned at intervals of $1/30$ th of a second, and the frame rate be 30 fps; then calculating one frame size would give us the number of bytes after which a RTS should be assigned in the incoming video data stream. The streams may then be sent to the playback sites where the media units in the audio and video stream with the same RTS are synchronized during playback.

Clearly, if data is compressed before being sent over the network the above strategy cannot be employed because compression algorithms such as JPEG are sensitive to the content of the image and the size of the compressed data may be variable. The requirement for assigning relative time-stamps to the various media streams hence calls for a mechanism based on the generation times of the data at the respective sites. A similar approach has been employed in [10]. However their approach may not be suitable if the application has some bounds on the accuracy with which the RTS should be assigned, because the accuracy with which their model can assign the RTS is bounded by the jitter of the network, which cannot be changed and may be large particularly in wide area networks.

In our approach for assigning the relative time-stamps, the data streams are time-stamped as they are generated at the

media sources, according to the respective clocks at the media sources and sent to the multimedia server. Note that time-stamp refers to the actual clock time when the media data is sent/received, while relative time-stamp is used to indicate an order among the various media units. By a mechanism which we will describe below, the various clock times are converted to a single clock time which we refer to as the *normalized clock time*. If the normalized clock time of the media units on the various media streams is the same, which indicates that they were generated at the same time, they are assigned the same RTS. It may be observed that this mechanism does not necessitate global clock synchronization. If the absolute time of all the clocks involved is known at some point in time and the clock skews of the different clocks is known, then it is possible to normalize the various clock times and hence determine the corresponding points in the different data streams that were generated at the same time by the media sources. This may be achieved as follows.

Consider the layout depicted in Figure 1. For the sake of simplicity, we consider only two active media sources in our discussion here, namely MS1 and MS2. A connection based on the quality of service (QOS) control model as described in [12] is established from the media sources (MS1 and MS2) to the multimedia server (MMS). The only services our protocol requires from the QOS session is a bound on the network jitter. Having established a QOS session, a packet called the *trigger packet* is sent from the MMS to the media sources to trigger them to send their respective clock times to the MMS. The trigger packet does not contain much information and thus constitutes little overhead. After an interval t another trigger packet is sent to the media sources on the same QOS session. On receiving the trigger packets, the media sources send a message back to the MMS indicating the time - according to their respective clocks - at which the two trigger packets were received.

Let x_0 and $x_0 + t$ be the times when the two trigger packets were received by MS1, according to MS1's clock and y_0 and $y_0 + t + p$ be the times when the two trigger packets were received by MS2 according to MS2's clock; for some p , where $-\infty < p < \infty$. If the delay incurred in sending the various packets were the same, then corresponding to a point in time x_1 according to the MS1 clock, the point in time according to the MS2 clock would be

$$y_x = \frac{(t+p) \times (x_1 - x_0)}{t} + y_0 \quad (1)$$

However since the delay is not constant for the different channels/packets, the maximum difference in delay that may be incurred by the different media streams is $2\delta_d$, where $\delta_d = \delta_{max} - \delta_{min}$ i.e. the jitter of the QOS session.

δ_{max} and δ_{min} are the maximum and the minimum delays respectively for the QOS session and are established at the time of setting up the QOS session. The maximum error ϵ that may be introduced due to the jitter, in determining the clock rates of the different sources is hence

$$\left(\left(\frac{(t+p+2\delta_d)(x_1-x_0)}{t} + y_0 \right) - \left(\frac{(t+p)(x_1-x_0)}{t} + y_0 \right) \right) = \frac{2\delta_d(x_1-x_0)}{t} \quad (2)$$

To achieve a lower value of ϵ would require that δ_d be small and t be large. Lower value of δ_d reflects the need for a higher quality of service parameter which would require greater communication resources to be dedicated. Having a greater value of t would require that the QOS session be maintained for a longer period of time. Since the session is required to have a high quality of service, blocking the communication resources for the time period t may be wasteful, since the only communication that is done over the QOS session in the interval t is the sending of two trigger packets which, as we mentioned, do not carry much data. An alternative is to terminate the QOS session soon after the first packet is sent and to reestablish a new QOS session for sending the second packet after the interval t . The maximum error in this case would be given by

$$\epsilon = \frac{(\delta_{d1} + \delta_{d2}) \times (x_1 - x_0)}{t} \quad (3)$$

where δ_{d1} is the jitter for the first QOS session and δ_{d2} is the jitter for the second session.

Having determined the absolute clock times of the different sources and the clock skew - after having received the reply of the two trigger packets - the QOS session is terminated. The regular connection for sending the media data from the various media sources to the MMS is then established and the data being sent from the different media sources is time-stamped according to the respective clocks and sent to the MMS. The MMS on receiving the media streams, assigns the same RTS to the media units on different media streams that have the same *normalized* time-stamp. Since the time-stamps are assigned at the media sources according to the respective clocks (which may be asynchronous), equation (1) is used to convert the different clock times to one clock time, which we refer to as the *normalized clock time*. The maximum error that may be introduced in comparing the different clock times is ϵ . Note that this value depends on the QOS parameters at the time of determining the clock times and the clock skew and is independent of the quality of the session during assigning of the relative time-stamps, and may hence be altered according to the requirement of the application.

4. Detection of Asynchronization

Once the data is stored in the multimedia server, the challenge is to play the various media streams at the destination, in a manner such that the media units with the same RTS are played synchronously. For the discussion in this section, we consider only simultaneous display synchronization. Other possible temporal constructs will be discussed in Section 6. Two factors playing the predominant role that may lead to asynchronization are:

- The network jitter i.e. $\Delta_{max} - \Delta_{min}$, where Δ_{max} is the maximum delay bound and Δ_{min} the minimum delay bound for the network connection while the media data is being sent. These values are different from δ_{max} and δ_{min} which are the maximum and the minimum delay bounds for the QOS session which precedes the actual transfer of media data.
- The different clock rates of the various media devices at the destination.

The steps followed for detecting asynchrony among the various media streams are similar to the ones used to assign the RTS to the media streams being emitted from the media sources. Initially a QOS session is established between the MMS and the various media destinations (MDs). Two trigger packets separated by time interval t are sent by the MMS to the media destinations. The MDs on receiving the packets note the time of arrival of the packets according to their respective clocks and send back this time information to the MMS. This information is used in equations (1) and (2) which are used to convert time from one clock scale to another and to estimate the maximum error that may have been introduced. The QOS session is then terminated.

In the next phase, after the QOS session has been terminated, the media data is sent to the various MDs by the MMS and the time at which the data is received at the MDs by the respective clocks is sent back to the MMS. Since the timings are obtained relative to different clocks, equation (1) is used to convert x_1 , measured by clock of MD1 to y_x which is the reference clock MD2. The maximum error that may have been introduced in the conversion is bounded by ϵ evaluated in equation (2).

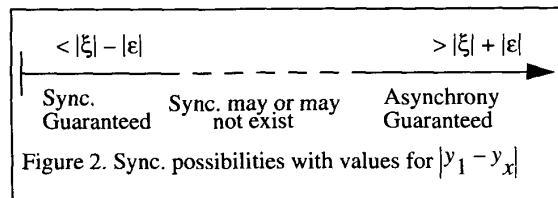
If it were possible to determine precisely the time at which the various media units with the same RTS on different media streams were displayed at the media destinations, it would enable us to determine with certainty if the media streams were synchronized or not. However, some error λ is introduced in evaluating equation (1) due to the difference in network delay. Moreover, this value λ is variable and depends on the network traffic conditions and may

vary from a minimum value 0 to the maximum value ϵ . Hence there would exist an interval for which it cannot be said with certainty if the media streams are synchronized or not.

Let the maximum permissible asynchrony between the media streams be ξ ; and y_1 and y_x be the times, according to the clock of MD2, of the media units from the two sources with the same RTS reaching the media destinations. We have the following cases:

- Synchronization is guaranteed if $|y_1 - y_x| \leq |\xi| - |\epsilon|$ (4)
- Asynchrony guaranteed if $|y_1 - y_x| > |\xi| + |\epsilon|$
- Synchronization may or may not exist otherwise.

Figure 2 illustrates the above conditions. The value of ξ is dependent on the application. Certain applications, for instance those involving audio, have tighter synchronization requirements and hence lesser ξ value than other applications, for instance those involving video. The desired value of ξ is provided by the application and ϵ can be computed from equation (3).



5. Re-synchronization

The policies for deciding when to trigger the synchronization mechanisms may be conservative, aggressive or probabilistic [2]. Conservative policies would trigger synchronization mechanisms only when $|y_1 - y_x| > |\xi| + |\epsilon|$, that is when asynchronization is guaranteed. Aggressive and probabilistic policies both trigger even while the value of $|y_1 - y_x|$ is in the dotted region of Figure 5, i.e. whenever there is any probability of the streams being out of synchronization. However the probabilistic policies base their decision on the statistical distribution of network delays and the playback periods, while the aggressive policies trigger the synchronization mechanisms whenever synchronization is not guaranteed.

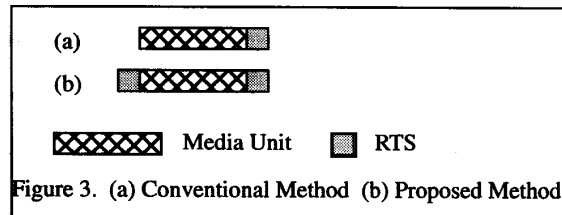
The time taken from the moment the asynchrony is detected to the time the streams are played back in synchronous fashion at the media destinations depends on the maximum network delay; since the asynchronization information has to travel from the media destinations to the

MMS, the necessary action taken by the MMS and then the data is sent to the media destinations. If the upper bound on the network delay is too large, the time taken may be unacceptable depending on the application. Hence an approach for predicting the time when the media streams could fall out of synchronization is desirable and has been suggested in the literature [10]. The accuracy with which the asynchrony may be *predicted* is dependent on the fractional deviation of the clocks from the nominal clock rate and the error bound in determining the asynchrony of the current media unit. Since our scheme can control the error bound ϵ , the accuracy of predicting the asynchrony may too be controlled.

6. Synchronization of Temporal Constructs

Synchronization refers to making events happen in a certain time order [4]. Most of the approaches dealing with synchronization in a distributed networked environment however, present synchronization in context of synchronization of all streams in parallel. In this section we describe how the mechanisms presented in earlier sections may be extended for synchronizing all the possible temporal constructs in an environment discussed in Section 2.

The mechanism for assigning relative time-stamps to media units of the different data streams is similar to the one proposed in Section 3. However the RTS are assigned both at the front and the rear end of the media units as shown in Figure 3. As is shown, if the RTS of a media unit is \mathcal{U} then \mathcal{U} is the assigned RTS at both the front and rear end of the media unit.



The process of detecting asynchrony, if any, among the media streams is in confluence to the discussion in Section 4. After the relationship between the clocks of the different media destinations is established, the media units P_a and P_b are sent on the (different) communication channels. At the media destination the times at which the front end and the back end RTS are received is noted, according to the respective clocks, and the information sent back to the MMS. The MMS on receiving this information and noting the type of synchronization desired

can, with the help of an extension of equation (4), described below, ascertain whether the media units are in synchronization. It may be noted here that the arrival time of both the front and rear end markers is noted and sent back, unlike the in-parallel synchronization scenario where only the time-stamp of the front marker is sent back.

It has been shown that there are thirteen distinct ways in which two time intervals may be related [5], [14]. Since six of them are inverse relations of the other [5], we discuss the details of the process of synchronization carried out by the MMS for only seven of these temporal relations.

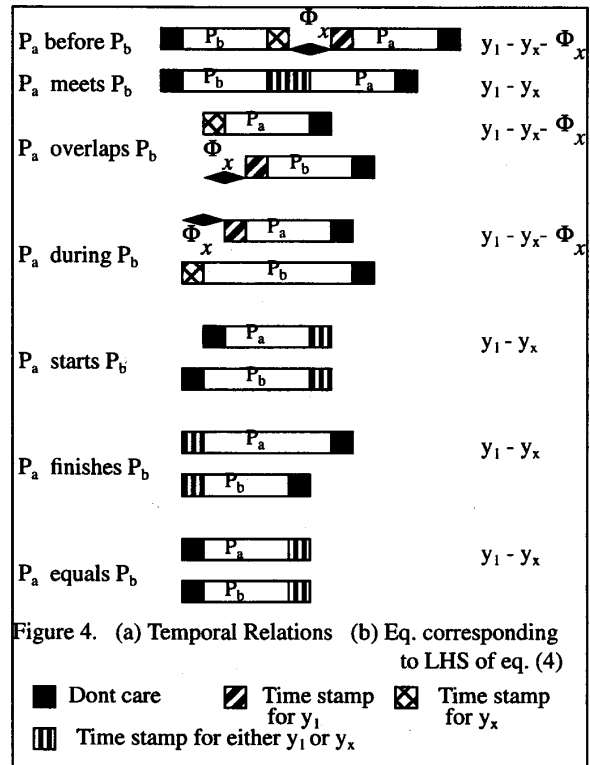


Figure 4 shows the seven different temporal relations and the corresponding equations that have to be used on the L.H.S. of equation (4). Consider the first temporal relation P_a before P_b . The requirement of this temporal relation is that P_a be rendered Φ_x time units before P_b . Since there is no time restriction as to when P_a begin or P_b end, the MMS needs to know the time-stamps (arrival time at the destination) of only the rear marker of P_a and the front marker of P_b . Since the front marker of P_a and the rear marker of P_b are not required for the determination of asynchrony they are marked *don't care*. For the equations (1) and (4), the time-stamp of the rear marker of P_a is substituted for the

value of y_1 and the time-stamp of the front marker of P_b is substituted for the value of y_x .

Substituting $|y_1 - y_x - \Phi_x|$ for $|y_1 - y_x|$ in equation (4), gives the equation for detecting the asynchrony for the *before* temporal relation, where y_1 is the time-stamp of the arrival of the rear RTS of the first media unit (P_a) and y_x is the time-stamp of the arrival of the front RTS of the second media unit, namely P_b .

That is if: (4a)

- $|y_1 - y_x - \Phi_x| \leq ||\xi| - |\epsilon||$, then the media streams are guaranteed to be in synchronization
- $|y_1 - y_x - \Phi_x| > ||\xi| - |\epsilon||$, then the media streams are assuredly out of synchronization
- Otherwise, it cannot be said with certainty if the streams are in synchronization or not.

Similarly the equations for all the other temporal relations are given in figure 4b). Re-synchronization of the other temporal constructs can be achieved in a manner similar to the one for synchronizing the parallel temporal construct.

In Figure 4, the *don't cares* refer to the markers whose time-stamp is not required by the MMS to detect asynchrony. As another example consider the temporal construct P_a *meets* P_b , the conditions for determining asynchronization remain the same as equation (4), with the only difference being the markers whose time-stamps are taken into consideration. We observe that for all the possible temporal relations only two forms of equations (4) exist, namely equation (4) itself and equation (4 a), with the difference among each of these equations being the markers that are considered for the substitution of y_1 and y_x .

7. Conclusions

We have proposed a scheme for synchronization of related media streams in a networked environment, where the accuracy of determining the asynchrony can be controlled, and hence tailored to the intended application. The idea of synchronizing media streams without global clock synchronization has been proposed in the literature, but the accuracy of detecting the asynchrony is fixed and hence their approach may not be suitable for certain applications particularly those in wide area networks environment. We have also proposed a mechanism where we ensure synchronization of all the possible temporal constructs and not just the in parallel construct, which is only one out of the thirteen possible temporal relations.

We are presently looking into mechanisms to express the temporal relationships while sending the data, so that differ-

ent temporal relationships may be interleaved and synchronized at the destination. Also we are simulating the protocol so that we may establish a mapping between the jitter guarantees provided by the network and the corresponding error bound with which the asynchrony may be detected.

References:

- [1] R. G. Herrtwich, "Time Capsules: An Abstraction for Access to Continuous-Media Data", IEEE Real-Time systems Symposium, Lake Buena Vista, 1990, IEEE press, Pages 11-20.
- [2] P. V. Rangan et. al, "Designing an On-Demand Multimedia Service", IEEE Communications Magazine, July 1992, Pages 56-64.
- [3] Sang H. Son and Nipun Agarwal, "Real-Time and Synchronization Issues in Multimedia Computing Systems", Workshop on the Role of Real-Time in Multimedia/Interactive Computing Systems, Nov 30, 1993, Raleigh, NC.
- [4] Ralf Steinmetz, "Synchronization Properties in Multimedia Systems", IEEE Journal on Selected Areas in Communication, Vol 8, No. 3, Pages 401-411, April 1990.
- [5] Thomas D. C. Little, Arif Ghafoor, "Synchronization and Storage Models for Multimedia Objects", IEEE Journal on Selected Areas in Communication, Vol. 8, No. 3, Pages 413-426, April 1990.
- [6] Petra Hoepner, "Synchronizing the Presentation of Multimedia Objects", IEEE Computer Communications, Vol. 15, No. 9, Nov 1992, Pages 557-564.
- [7] David P. Anderson, George Homsy, "A Continuous Media I/O Server and Its Synchronization Mechanism", IEEE Computer, Pages 51-57, October 1991.
- [8] C. Nicolaou, "An Architecture for Real-Time Multimedia Communication Systems", IEEE Journal on Selected Areas in Communication, Pages 391-400, Vol. 8, No. 3, April 1990.
- [9] Julio Escobar, D. Deutsch, C. Partridge, "Flow Synchronization Protocol", Globecom 1992.
- [10] P. V. Rangan et. al, "Techniques for Multimedia Synchronization in Network File Systems", IEEE Computer Communications, Vol. 16, No. 3, March 1993, Pages 168-176.
- [11] Sang H. Son and Nipun Agarwal, "Synchronization of Temporal Constructs in Distributed Multimedia Systems with controlled accuracy", Tech Rep. # CS-93-57, Dept. of Computer Science, University of Virginia, Oct 28, 1993.
- [12] Hideyuki Tokuda et. al., "Continuous Media Communication with Dynamic QOS Control Using ARTS with an FDDI Network", SIGCOMM 1992, Baltimore, MD, Pages 88-98, August 1992.
- [13] Ralf Steinmetz, "Multimedia Synchronization Techniques: Experience Based on Different System Structures", IEEE Multimedia Workshop 1992, Monterey, April 1992.
- [14] C. L. Hamblin, "Instants and Intervals", Proc. of 1st International Soc. for the Study of Time, J.T.Fraser et. al, Eds., New York: Springer-Verlag, 1972, Pages 324-331.
- [15] D. Ferrari et. al., "A Scheme for Real-Time Channel Establishment in Wide-Area Networks", IEEE journal on Selected Areas in Communication, Vol. 7, No. 3, April 1990, Pages 368-379.